How Compositional Generalization and Creativity Improve as Diffusion Models are Trained

Alessandro Favero* EPFL, Lausanne, Switzerland

Antonio Sclocchi* UCL, London, United Kingdom

Francesco Cagnetta SISSA, Trieste, Italy

Pascal Frossard EPFL, Lausanne, Switzerland

Matthieu Wyart

Johns Hopkins University, Baltimore, U.S. & EPFL, Lausanne, Switzerland

ALESSANDRO.FAVERO@EPFL.CH

A.SCLOCCHI@UCL.AC.UK

FRANCESCO.CAGNETTA@SISSA.IT

PASCAL.FROSSARD@EPFL.CH

MWYART1@JH.EDU

Abstract

Natural data is often organized as a hierarchical composition of features. How many samples do generative models need to learn the composition rules, so as to produce a combinatorially large number of novel data? What signal in the data is exploited to learn those rules? We investigate these questions in the context of diffusion models both theoretically and empirically. Theoretically, we consider simple probabilistic context-free grammars—tree-like graphical models used to represent the hierarchical and compositional structure of data such as language and images. We demonstrate that diffusion models learn the grammar's composition rules with the sample complexity required for clustering features with statistically similar context. This clustering emerges hierarchically: higher-level features associated with longer contexts require more data to be identified. This mechanism leads to a sample complexity that scales polynomially with the said context size. As a result, diffusion models trained on an intermediate dataset size generate data coherent up to a certain scale, but that lacks global coherence. We test these predictions in different domains and find remarkable agreement: both generated texts and images achieve progressively larger coherence lengths as the training time or dataset size grows.

1. Introduction

Compositional generalization, the ability to understand and generate novel combinations of known components, is a fundamental characteristic of human intelligence and creativity. For instance, this skill allows humans to create grammatically correct and meaningful sentences never heard before or to reason originally by assembling together known ideas. Under which conditions can machines learn such a skill? The success of diffusion models in producing realistic data across various domains [3, 19, 34, 43, 44] provides a unique opportunity to study how this ability emerges. Fundamental questions include: What signals in the data are exploited by neural networks to learn the

^{*} Equal contribution.

compositional rules? How many training examples are needed to learn such rules, and in what order are they learned? How does the finiteness of the training set affect the structure of generated data?

To address these questions theoretically, we bridge two viewpoints developed in the context of natural language processing. On the one hand, *symbolic approaches* aim to describe the structure of data via a list of rules that generate them. For example, probabilistic *context-free grammars* (PCFG) [10] describe sentences with trees, whose nodes are hidden variables that can generate other nodes or leaves according to probabilistic production rules. PCFGs can approximate both structural and semantic aspects of text and have been proposed for the description of images under the name of *Pattern Theory* [16, 22, 42]. On the other hand, *statistical approaches* use data-driven analyses agnostic to expert knowledge of grammatical structure. A notable example is *word2vec* [29], where a shallow network learns meaningful representations of words by predicting their neighborhood.

We unify these two viewpoints by studying how diffusion models learn simple PCFGs. In particular, we show empirically that the learning process of diffusion models is hierarchical, progressively capturing compositional rules at deeper levels of the PCFG's hierarchy. We then argue that the grammar rules can be deduced iteratively by clustering, as in word2vec, sequences of tokens based on the statistics of their context. For each level, we analytically derive the corresponding sample complexity. We show that it matches the number of data required by the diffusion model to generate data that follow the PCFG rules up to that level. Since this hierarchical clustering procedure requires a number of samples that is polynomial in the size of the token's sequence, this mechanism allows the diffusion model to learn a high-dimensional distribution while avoiding the *curse of dimensionality*. Beyond simple PCFGs, we predict that diffusion models trained on limited samples generate data that is locally coherent (i.e., satisfying local compositional rules), but not globally, with a coherence length growing with the training time/number of samples. We confirm this prediction in diffusion models trained on OpenWebText and ImageNet.

2. Setup

Diffusion models Denoising diffusion models are a family of generative models built to draw samples from a target distribution by reversing a process where noise is gradually added [19, 43–45]. Let t denote the time index running in [0, ..., T], and let $q(\cdot)$ be the distribution we aim to sample from, with $x(0) \sim q(x(0))$ denoting a sample from this distribution. A diffusion model is composed of two main parts. Firstly, a *forward process* that sequentially adds noise to the data to produce the sequence $\{x(t)\}_{1 \le t \le T}$, $q(x(1), ..., x(T) \mid x(0)) = \prod_{t=1}^{T} q(x(t) \mid x(t-1))$, culminating in a purely noisy sample x(T). Secondly, a *backward process* that reverses the noise addition step by step and is typically learned by training a neural network to approximate the backward transition kernels $p(x(t-1) \mid x(t))$. This process effectively learns the *score function*, which is proportional to the conditional expectation $\mathbb{E}_{q(x(0)|x(t))}[x(0)] \coloneqq \mathbb{E}[x(0)|x(t)]$. To draw a new sample from $q(\cdot)$, one starts with a noise sample $x(T) \sim q(x(T))$ and then applies the learned backward process to obtain a clean sample $x(0) \sim q(x(0))$. Depending on the characteristics of the data space, diffusion models differ in how they define the forward process.

Probabilistic graphical models To investigate how diffusion models learn compositional structures, we consider synthetic datasets generated via a *probabilistic context-free grammar* (PCFG) [36]: a collection of symbols and rules that prescribe how to generate sequence data starting from a single feature. Generic PCFGs consist of a vocabulary of hidden (*nonterminal*) symbols, a vocabulary of visible (*terminal*) symbols and *production rules* that quantify the probability that one hidden symbol



Figure 1: Learning different levels of the grammar. (a) Accuracy at various levels as a function of training set size P. (b) Online learning setting, where fresh training points are sampled at each step. (c) Token-token correlations for $N = 10^6$ samples generated by the diffusion model trained with P training points.

generates tuples of either hidden or visible symbols. The Random Hierarchy Model (RHM) [8] is a particular PCFG, including the following additional assumptions to make it analytically tractable. (i) The nonterminal symbols are split into L finite vocabularies $(\mathcal{V}_{\ell})_{\ell=1,...,L}$ of finite size v and $\mathcal{V} \equiv \mathcal{V}_0$ denotes the vocabulary of terminal symbols. (ii) All the production rules transform one level- $(\ell + 1)$ symbol into a string of s level- ℓ symbols. (iii) There are m unambiguous production rules per nonterminal symbol, i.e., two distinct nonterminals cannot generate the same s-tuple. The rules are randomly chosen and frozen for a given instance of the RHM. We call the m strings produced by any given symbol synonyms; (iv) All the available production rules are equally likely. Due to assumptions (i) and (ii), the data-generating process can be represented as a regular tree graph with depth L and branching ratio s. The leaf nodes ($\ell = 0$) correspond to the tokens of the visible data, which form strings of size $d = s^L$. Knowing the production rules and the tree structure of the RHM, in these models the score can be computed exactly using Belief Propagation. In what follows, we study instead how many samples are necessary to learn sampling from the distribution and the score function from data.

3. How diffusion models learn a grammar

In this section, we investigate how diffusion models learn to generate data from the RHM. Given an instantiation of the RHM, we uniformly sample P distinct training points, i.e., strings from the grammar. Each input symbol is encoded as a *one-hot vector*, $x \in \{0,1\}^{d,v}$. With this dataset, we train a *Discrete Denoising Diffusion Probabilistic Model* (D3PM) [2] with uniform transition probabilities [20]. The diffusion model architecture is a convolutional U-Net [35] with L resolution blocks in both the encoder and decoder. We use the neural network to predict the conditional expectation $\mathbb{E}(x(0)|x(t))$, which parameterizes the reverse diffusion process. We explore both an offline learning setting, where a finite dataset is generated, and the model is trained over multiple epochs, and an online learning setting. Additional details are reported in Appendix E.

After training, we generate 1024 samples and evaluate whether the generated data satisfies the compositional rules of the RHM at different hierarchical levels. Specifically, we define the *accuracy* A_{ℓ} at level ℓ as the fraction of generated samples that satisfy level- ℓ rules. Figure 1a shows the accuracy at different levels as a function of P. The results reveal a staged learning process: the low-

level rules, governing local structures, are learned first, followed by progressively higher-level rules that enforce global coherence. Thus, models trained on intermediate P values generate data that are locally consistent but lack global coherence. The inset of Figure 1a compares favorably the scaling of accuracy with our theoretical prediction, which is derived in the next section. This prediction indicates that learning to satisfy rules at level ℓ requires a number of samples that scales as $m^{\ell+1}$. Importantly, this scaling is polynomial, not exponential, in the data dimension $d = s^L$. Specifically, the sample complexity to learn all rules is $m^{L+1} = md^{\log m/\log s}$. Figure 1b demonstrates that the same staged learning process applies in the online learning setting, where fresh training samples are drawn at each training step. This progressive acquisition of compositional rules also appears in the internal correlations of the generated sequences, defined as the Frobenius norm of the covariance matrix between two visible tokens at distance t. As shown in Figure 1c, at small training set sizes or training times, only nearby tokens exhibit significant correlations, while long-range correlations approach sampling noise (black dashed line, given by $1/(vN^{1/2})$, where N is the number of sequences used to measure correlations). As training progresses, long-range correlations emerge.

Emergence of hierarchical representations To generate sequences that satisfy the compositional rules of the RHM, the diffusion model needs to construct internal representations of the latent variables at each level. To do so, it must represent together inputs that differ by low-level synonyms. In Appendix F, we show that this is indeed the case: as the training set size increases, the hidden representations of the U-Net become insensitive to higher and higher levels of synonyms.

4. Theoretical Analysis

Let $\mathbf{x}(t)$ denote a noised RHM string obtained from a clean sample $\mathbf{x}(0)$. Diffusion models are trained to predict the conditional expectation $\mathbb{E}[\mathbf{x}(0) \mid \mathbf{x}(t)]$. We analyze the low-noise limit $(t \to 0)$, assuming that only the first token $x_1(t)$ is corrupted. In this case, the network must therefore estimate $\mathbb{E}[x_1(0) \mid \mathbf{x}_{2:d}(0)]$, which is proportional to the correlation between $x_1(0)$ and the context $\mathbf{x}_{2:d}(0)$. In the RHM, these correlations are invariant under exchanges of synonyms [7]. Hence, the score depends only on the latent variables that generated the context:

$$\mathbb{E}[x_1(0) \mid \boldsymbol{x}_{2:d}(0)] = \mathbb{E}[x_1 \mid \boldsymbol{x}_{2:s}, \mathbf{h}_{2:s}^{(1)}, \mathbf{h}_{2:s}^{(2)}, \dots, \mathbf{h}_{2:s}^{(L-1)}],$$
(1)

where $h_i^{(\ell)}$ is a latent at level ℓ . This expression depends only on a sequence of length $(s-1)L \ll d = s^L$. Thus, learning the production rules and clustering synonyms by their latents drastically reduces the effective dimensionality of the task.

In Appendix B, we compute the sample complexity needed to reconstruct latent variables. By equating the magnitude of their correlations with the visible tokens with the sampling noise in a finite dataset, we find that the number of samples required to reconstruct level $\ell - 1$ latents scales as

$$P^{(\ell)} \sim m^{\ell+1}.\tag{2}$$

As a result, lower-level production rules are learned first, while deeper rules require progressively more data. This scaling exactly matches the U-Net learning curves in Figure 1.

5. Natural data

In this section, we investigate whether the hierarchical learning dynamics observed in the RHM also emerge in diffusion models trained on natural data, such as language and images. Since both

10⁸ training tokens

In popular spokesman typeted in diversity adventure allow price Zha Tampa usually Pages superstays's under leveldowns swim a cycle who retains highly weapons batch floor despite

10⁹ training tokens

Just like you are growing fast and growing strong. But this way you became organic, changed someone else 2019s. But even then you made them off. I sort came to smile around, because I was in China okay.

10¹⁰ training tokens

At the beginning of winter when I walked around; even if he would be talking to me, on the highest field and back in the second round in my team I would take him over in his cell because it was my game against Juventus. MD4 on OpenWebText G_{1}^{0} G_{2}^{0} G_{2}^{0

(a) Text generated at different training stages.

(b) Correlations of generated text.

Figure 2: Stage-wise learning of masked language diffusion model on OpenWebText. (a) The text generated by MD4 exhibits longer coherence spans as training progresses. (b) Correlations between tokens at a distance t in the generated text.

modalities have an inherent compositional structure—where words form sentences and object parts form images—we expect their learning process to progress hierarchically.

Language diffusion models We consider MD4 [41], a state-of-the-art masked diffusion model with absorbing state for discrete data such as language, as described in Appendix E. We train MD4 from scratch using a standard GPT-like transformer architecture with 12 layers ($\approx 165M$ parameters) on the OpenWebText corpus [15]. The model is trained for a full epoch on the training split ($\approx 10^{10}$ tokens) using the same hyperparameters as Shi et al. [41]. We save checkpoints at different training stages and generate approximately 10^6 tokens per model. Figure 2a presents text samples generated at various training times. Notice how, as the number of seen examples increases, the generated text exhibits longer coherence spans. In particular, the intermediate checkpoint ($\approx 10^9$ tokens) correctly assembles words locally but fails to generate coherent sentences, similar to what we observed in our synthetic experiments in Section 3. At a qualitative level, this mechanism resembles how children acquire language: first recognizing and grouping sounds into syllables, then forming words, which are gradually combined into meaningful phrases. We confirm this result quantitatively by measuring the token-token correlation function of the generated text (Figure 2b), as done for the RHM. Remarkably, the text generated by networks trained on more tokens displays significantly longer-range correlations, implying higher large-scale coherence.

Vision diffusion models For image data, we consider Improved *Denoising Diffusion Probabilistic Models* (DDPMs) [30]. Specifically, we train a U-Net model architecture [35, 37] with multi-head attention layers [47] ($\approx 120M$ parameters). The model is trained for 10 epochs on ImageNet 64×64 . Figure 8a in Appendix F, illustrates images generated at different training stages. Initially, the outputs exhibit patterns of textures. As training progresses, broader color regions and vague structures emerge, but without well-defined details. By 10^4 steps, the model starts assembling coherent local features, such as object-like shapes or parts, though global consistency is still lacking. Finally, images from the last checkpoint exhibit highly structured and realistic compositions, indicating that the model successfully learns to generate coherent scenes with well-defined objects. To quantify these observations, in Appendix F we analyze the hierarchical and compositional structure of generated images using deep latent representations from a pre-trained convolutional architecture, showing that structures in deeper layers emerge later in training.

6. Conclusions

We have provided a theory explaining how diffusion models can learn certain structured distributions with a polynomial number of data in the dimension, thus beating the curse of dimensionality. We showed that if data consists of a hierarchical combination of features, U-Nets can lower the data dimension by giving identical representations to groups of features that have similar contexts. This idea, explicit in word2vec, is performed hierarchically in diffusion models. This framework predicts that as the training time or training set size increases, generated data becomes coherent at larger scales. We provided direct evidence that this is the case for generated text and images.

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv preprint arXiv:2305.13673*, 2023.
- [2] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- [3] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3), 2023.
- [4] Giulio Biroli and Marc Mézard. Generative diffusion in very large dimensions. Journal of Statistical Mechanics: Theory and Experiment, 2023(9):093402, 2023.
- [5] Adam Block, Youssef Mroueh, and Alexander Rakhlin. Generative modeling with denoising auto-encoders and langevin sampling. *arXiv preprint arXiv:2002.00107*, 2020.
- [6] Francesco Cagnetta and Matthieu Wyart. Towards a theory of how the structure of language is acquired by deep neural networks. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [7] Francesco Cagnetta, Alessandro Favero, and Matthieu Wyart. What can be learnt with wide convolutional neural networks? In *International Conference on Machine Learning*, pages 3347–3379. PMLR, 2023.
- [8] Francesco Cagnetta, Leonardo Petrini, Umberto M Tomasini, Alessandro Favero, and Matthieu Wyart. How deep neural networks learn compositional data: The random hierarchy model. *Physical Review X*, 14(3):031001, 2024.
- [9] Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. arXiv preprint arXiv:2302.07194, 2023.

- [10] Noam Chomsky. Aspects of the Theory of Syntax. Number 11. MIT press, 2014.
- [11] Hugo Cui, Florent Krzakala, Eric Vanden-Eijnden, and Lenka Zdeborová. Analysis of learning a flow-based generative model from limited sample complexity. *arXiv preprint arXiv:2310.03575*, 2023.
- [12] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. *Proceedings of Thirty-Fifth Conference on Learning Theory*, 178:5413–5452, 2022.
- [13] Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. arXiv preprint arXiv:2208.05314, 2022.
- [14] Jérôme Garnier-Brun, Marc Mézard, Emanuele Moscato, and Luca Saglietti. How transformers learn structured data: insights from hierarchical filtering. arXiv preprint arXiv:2408.15138, 2024.
- [15] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007. github.io/OpenWebTextCorpus, 2019.
- [16] Ulf Grenander. *Elements of pattern theory*. JHU Press, 1996.
- [17] Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex Smola. A kernel method for the two-sample-problem. *Advances in neural information processing systems*, 19, 2006.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, June 2016. doi: 10.1109/CVPR.2016.90.
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [20] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- [21] Sadeep Jayasumana, Srikumar Ramalingam, Andreas Veit, Daniel Glasner, Ayan Chakrabarti, and Sanjiv Kumar. Rethinking fid: Towards a better evaluation metric for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9307–9315, 2024.
- [22] Ya Jin and Stuart Geman. Context and hierarchy in a probabilistic image model. In 2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06), volume 2, pages 2145–2152. IEEE, 2006.
- [23] Zahra Kadkhodaie, Florentin Guth, Stéphane Mallat, and Eero P Simoncelli. Learning multiscale local conditional probability models of images. arXiv preprint arXiv:2303.02984, 2023.
- [24] Mason Kamb and Surya Ganguli. An analytic theory of creativity in convolutional diffusion models. arXiv preprint arXiv:2412.20292, 2024.

- [25] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.
- [26] Eran Malach and Shai Shalev-Shwartz. A provably correct algorithm for deep learning that actually works. arXiv preprint arXiv:1803.09522, 2018.
- [27] Song Mei. U-nets as belief propagation: Efficient classification, denoising, and diffusion in generative hierarchical models. arXiv preprint arXiv:2404.18444, 2024.
- [28] Song Mei and Yuchen Wu. Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models. arXiv preprint arXiv:2309.11420, 2023.
- [29] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information* processing systems, 26, 2013.
- [30] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- [31] Maya Okawa, Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. Compositional abilities emerge multiplicatively: Exploring diffusion models on a synthetic task. Advances in Neural Information Processing Systems, 36, 2024.
- [32] Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution estimators. arXiv preprint arXiv:2303.01861, 2023.
- [33] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017. doi: 10.23915/distill.00007. https://distill.pub/2017/feature-visualization.
- [34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pages 10684–10695, 2022.
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9,* 2015, proceedings, part III 18, pages 234–241. Springer, 2015.
- [36] Grzegorz Rozenberg and Arto Salomaa. Handbook of Formal Languages. Springer, January 1997. doi: 10.1007/978-3-642-59126-6.
- [37] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. arXiv preprint arXiv:1701.05517, 2017.
- [38] Antonio Sclocchi, Alessandro Favero, Noam Itzhak Levi, and Matthieu Wyart. Probing the latent hierarchical structure of data via diffusion models. *arXiv preprint arXiv:2410.13770*, 2024.

- [39] Antonio Sclocchi, Alessandro Favero, and Matthieu Wyart. A phase transition in diffusion models reveals the hierarchical nature of data. *Proceedings of the National Academy of Sciences*, 122(1):e2408799121, 2025.
- [40] Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the ddpm objective. arXiv preprint arXiv:2307.01178, 2023.
- [41] Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis K Titsias. Simplified and generalized masked diffusion for discrete data. arXiv preprint arXiv:2406.04329, 2024.
- [42] Jeffrey Mark Siskind, J Sherman, Ilya Pollak, Mary P Harper, and Charles A Bouman. Spatial random tree grammars for modeling hierarchal structure in images with regions of arbitrary shape. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(9):1504–1519, 2007.
- [43] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- [44] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- [45] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020.
- [46] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.
- [48] Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. *arXiv preprint arXiv:2011.14522*, 2020.
- [49] Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Rewarddirected conditional diffusion: Provable distribution estimation and reward improvement. arXiv preprint arXiv:2307.07055, 2023.

Appendix A. Related work

Sample complexity in diffusion models Under mild assumptions on the data distribution, diffusion models exhibit a sample complexity that scales exponentially with the data dimension [5, 32]. It is not the case if data lie on a low-dimensional latent subspace [9, 13, 49], correspond to Gaussian mixture models [4, 11, 40], Ising models [28], or distributions that can be factorized across spatial scales [23]. These works do not consider the sample complexity in compositional data.

Compositional generalization of diffusion models Okawa et al. [31] considered synthetic compositional data to empirically show how diffusion models learn to generalize by composing different concepts, in the absence of a compositional hierarchy. Kamb and Ganguli [24] studied how equivariant diffusion models can compose images by combining local patches seen in the dataset. Sclocchi et al. [38, 39] showed that diffusion on hierarchically compositional data can be solved using Belief Propagation. Mei [27] showed that U-Nets can efficiently approximate the Belief Propagation algorithm on hierarchical data. Yet, efficient representability does not guarantee learnability by gradient descent for hierarchical data [8]. These works do not, however, address the sample complexity of diffusion models learned by gradient descent or variations of it.

Learning hierarchical representation via next-token prediction It has been observed that transformers trained on next-token prediction on PCFGs learn a hierarchical representation of the data that reflects the structure of the latent variables [1, 6, 14]. Closest to our work, Cagnetta and Wyart [6] showed that for the prediction of the last token in a sequence of fixed length, the latent structure is learned hierarchically, with a sample complexity polynomial in the context length. Our work extends this finding to diffusion models, in a setup where complete sequences can be generated. This setup allows us to make novel predictions on the limitations of generated data as a function of the training set size, which we test empirically across domains.

Appendix B. Theoretical analysis

To derive the sample complexity of the U-Net, we build upon prior work that explains how deep networks efficiently learn hierarchical tasks. This result is achieved by building a lower-dimensional representation that iteratively clusters synonyms [26], allowing the network to recover the latent hierarchical structure of the data. This clustering mechanism is based on statistical correlations between *s*-tuples of tokens and the given task—supervised or self-supervised—which are identical for synonyms. Notably, the sample complexity of deep networks trained with gradient descent aligns with the training set size required to detect these correlations [6, 8]. For supervised learning, this connection can be justified in a one-step gradient descent (GD) setting.

Here, we extend these results to diffusion models. First, we demonstrate that learning the score function in the low-noise limit corresponds to a task invariant to exchanging synonyms, and could thus be simplified by reconstructing the latent variables. Then, we compute the sample complexities required to reconstruct latent variables of different levels using correlations. We conclude by showing that a a clustering algorithm based on correlations does indeed recover the latent variables with the predicted sample complexities and b) the sample complexity required to reconstruct first-level latent variables can be recovered in a one-step-GD setting.



(a) U-Net scheme.

(b) RHM structure.

Figure 3: U-Net scheme and RHM structure. (a) To denoise the RHM data, the U-Net has to predict the conditional expectation $\mathbb{E}[\boldsymbol{x}(0)|\boldsymbol{x}(t)]$ for a given noisy input $\boldsymbol{x}(t)$, which is proportional to the correlations of the single tokens $x_i(0)$ with $\boldsymbol{x}(t)$. This can be done efficiently by learning the latent hierarchical structure of the data. (b) The correlations of the RHM data reflect the tree structure of the model (black squares represent the rules at different levels). For the token x_1 , using the correlations with tuples at different levels (highlighted in red), the conditional expectation $\mathbb{E}[x_1|\boldsymbol{x}_{2:8}]$ can be represented as $\mathbb{E}[x_1|\boldsymbol{x}_2, h_2^{(1)}, h_2^{(2)}]$.

B.1. Learning the score in the low-noise limit

Input-output correlations in diffusion models The loss function of diffusion models is minimized when the model prediction converges to the conditional expectation $\mathbb{E}[\boldsymbol{x}(0)|\boldsymbol{x}(t)]$, which is sampled in the limit of infinite diffusion trajectories and is proportional to the score function [2, 43, 44]. Since the expectation operates independently for each v-dimensional one-hot-encoded token $x_j(0), j \in [d]$, we have that $\mathbb{E}[x_j(0)|\boldsymbol{x}(t)]$ is directly proportional to the correlation between a token $x_j(0)$ and the input $\boldsymbol{x}(t)$.

Score function at low noise We now consider a small-noise regime $t \to 0$ where only the first token has been changed by noise, to some value $x_1(t)$ uncorrelated with $x_1(0)$. In this case, the function that the network has to learn is $\mathbb{E}[x_1(0)|\mathbf{x}_{2:d}(0)]$, proportional to the correlations of the first token with the remaining sequence of length d-1. Since these correlations are invariant under exchanges of synonyms [8], they correspond to the correlations of the x_1 token with the latents at all levels generating the rest of the sequence, i.e., $\mathbb{E}[x_1|\mathbf{x}_{2:s}, \mathbf{h}_{2:s}^{(1)}, \mathbf{h}_{2:s}^{(2)}, \dots, \mathbf{h}_{2:s}^{(L-1)}]$ (Figure 3(b)). This function depends on a sequence of length (s-1)L, much smaller than the data dimension $d = s^L$. In other words, knowing the latent variables allows for a significant reduction of the problem dimensionality.

B.2. Sample complexities

In this section, we determine the sample complexities required to reconstruct the tuple of latent variables of different levels $\mathbf{h}_{2:s}^{(\ell)}$ appearing in the low-noise score function. As shown in Cagnetta and Wyart [6], latents can be reconstructed via their correlations with the noised token x_1 . We thus work under the following assumption.

Assumption The U-Net learns to generate data that is consistent with the rules at layer ℓ when the correlations between a visible token and a tuple of latents at layer $\ell - 2$ become detectable from the training data.

Hence, in what follows, we compute the number of samples required to detect these correlations.

Local constraints The first step in the learning process is to recognize the valid *s*-tuples generated by the RHM at the visible level. Since these tuples lack internal structure, they can only be memorized. Each tuple can take vm possible configurations corresponding to v symbols for the first-level latents and m representations (synonyms) for each of them. Thus, the sample complexity required to learn the local constraints scales as $P^{(1)} \sim vm$.

First-level latents Once the local constraints are learned, the network can refine its estimate of x_1 by utilizing correlations with the neighboring tuples $x_{s+1:2s}, \ldots, x_{s^2-(s-1):s^2}$. The sample complexity required to detect the correlations between x_1 and $x_{s+1:2s}$ was computed in Cagnetta and Wyart [6] and corresponds to $P_{corr}^{(2)} \sim vm^3$. After learning the first-level rules, the network can collapse the $(s^2 - s)$ -dimensional sequence of neighboring tuples into the corresponding first-level latents $\mathbf{h}_{2:s}^{(1)}$.

Second-level latents Having built the first-level latent representation, the model can leverage correlations between *s*-tuples of first-level latents $h_i^{(1)}$'s and the first token to learn the rules at the second level, further improving the denoising task. These correlations can be computed by studying the statistics of the token-latent tuple correlations,

$$\mathbb{P}[x_1 = \mu, \mathbf{h}_{s+1:2s}^{(1)} = \nu] - \mathbb{P}[x_1 = \mu] \mathbb{P}[\mathbf{h}_{s+1:2s}^{(1)} = \nu],$$
(3)

over realizations of the RHM. Since correlations have zero mean, we take the standard deviation over RHM realizations as an estimate of their typical size. As shown in Appendix C, the resulting correlation magnitude is given, in the limit of large v and m, by $C^{(3)} \simeq (v^3 m^5)^{-1/2}$. Since a finite training set of size P only allows measuring the empirical correlation function, we compare the magnitude of correlations with the sampling noise, which has magnitude $(v^2 m P)^{-1/2}$. Thus, the number of samples required to detect correlations between tuples of first-level latents and visible tokens, denoted as $P_{\text{corr}}^{(3)} \sim vm^4$.

Extension to general depth ℓ The same procedure generalizes to any depth ℓ . The correlations between tuples of latents at level $\ell - 2$ and visible tokens, having lowest common ancestor at level ℓ , have magnitude $C^{(\ell)} \simeq \sqrt{1/(v^3m^{\ell+2})}$. Meanwhile, the sampling noise remains of order $(v^2mP)^{-1/2}$. Equating these terms gives the sample complexity required to reconstruct level- $(\ell - 1)$ latents,

$$P_{\rm corr}^{(\ell)} \sim v m^{\ell+1}.\tag{4}$$

This result indicates that learning rules leveraging correlations at depth L requires a number of samples scaling as $m^{L+1} = md^{\log m/\log s}$, which is polynomial (and not exponential) in the dimension. Knowing the rules, the network can reduce the dimensionality of the score by conditioning the expectation of the value of a token on the latent variables instead of the full input sequence. Remarkably, Eq. (4) displays the same scaling observed in our experiments with the U-Net in Section 3, confirming our assumption.

B.3. Clustering and one-step GD

Clustering To validate the hypothesis that synonyms can be grouped based on correlations, we consider a simple clustering algorithm that computes the empirical correlations between (latent) tuples and a visible token and then applies k-means clustering. As shown in Figure 4, the sample complexity for such an algorithm (black points) closely follows the theoretical prediction $P_{\text{corr}}^{(L)} \sim m^{L+1}$. We also test a modified algorithm that uses all the tokens in the first visible tuple instead of just the first (red points in Figure 4). Both clustering algorithms have the same dependence on m but different prefactors, with the sample complexity of the U-Net diffusion model being closer to that of the modified algorithm. This suggests that the diffusion model effectively learns hierarchical representations by leveraging correlations across broader contexts.

One-step gradient descent Finally, to support the connection with standard training techniques, we consider a simplified setting where a linear architecture is trained via gradient descent to predict the token x_{s+1} given an adjacent tuple $(x_1, \ldots x_s)$. This task corresponds to learning the score function $\mathbb{E}[x_{s+1}(0)|\mathbf{x}_{1:s}(0)]$, which is invariant to exchanging the tuple $(x_1, \ldots x_s)$ with a synonym. As proved in Appendix D, one step of gradient descent aligns the learned weights with the empirical token-tuple correlations. Consequently, if the size of the training set is large enough for the accurate measure of correlations, then the network can build a representation of the tuple $(x_1, \ldots x_s)$, which is invariant to exchanging synonyms. This invariance is empirically observed for the U-Net in Figure 5 of Appendix F.

Appendix C. Token-latent tuple correlations

In this section, we derive our estimate for the magnitude of the correlations between x_1 and tuples of latent, level- $(\ell - 1)$ features $h_{(i-1)\times s+1:i\times s}^{(\ell-1)}$, with $i = 2, \ldots, s$ and $\ell = 1, \ldots, L-1$ (level-0 latents $h^{(0)}$ correspond to visible tokens). These correlations are identical for all the tuples of latents corresponding to the same higher-level feature $h_i^{(\ell)}$, thus can be used to reconstruct level- ℓ latents. For instance, with s = 2, so that i = 2 (see Figure 3), the correlations of x_1 with (x_3, x_4) determine the value of $h_2^{(1)}$, while those with $(h_3^{(1)}, h_4^{(1)})$ determine $h_2^{(2)}$. To simplify the notation, we will stick to the case i = 2 for the remainder of the section. Then, the goal is to compute the statistics of

$$C^{(\ell+1)}(\mu,\boldsymbol{\nu}) \coloneqq \mathbb{P}\left\{X_1 = \mu, \mathbf{h}_{s+1:2s}^{(\ell-1)} = \boldsymbol{\nu}\right\} - \mathbb{P}\left\{X_1 = \mu\right\} \mathbb{P}\left\{\mathbf{h}_{s+1:2s}^{(\ell-1)} = \boldsymbol{\nu}\right\},\tag{5}$$

over realizations of the RHM.

For each visible token i = 1, ..., d, single-token probabilities can be written as products of probabilities over the single production rules,

$$\mathbb{P}\left\{X_{i}=\mu\right\} = \sum_{\mu_{1},\dots,\mu_{L}=1}^{\nu} p_{i_{1}}^{(1)}(\mu|\mu_{1})\dots p_{i_{L}}^{(L)}(\mu_{L-1}|\mu_{L})p^{(L+1)}(\mu_{L}),\tag{6}$$

where

- (i) the indices i_L, \ldots, i_L are such that $i_L \ldots i_1$ equals the *s*-ary representation of *i*, with $i_\ell = 1, \ldots, s$, and 1's added to ensure that the representation always consists of *L* indices. In other words, the multi-index i_L, \ldots, i_L uniquely identifies the path linking the root of the tree to the *i*-th leaf.
- (ii) $p_{i_{\ell}}^{(\ell)}(\mu_{\ell-1}|\mu_{\ell})$ denotes the probability of choosing, among the available production rules starting from μ_{ℓ} , one that has the symbol $\mu_{\ell-1}$ on the i_{ℓ} -th position of the right-hand size.
- (iii) $p^{(L)}(\mu_L)$ denotes the probability of selecting the symbol μ_L as the root (1/v for our model).

These decompositions arise naturally due to the connection between probabilistic context-free grammars and Markov processes. Similar decompositions apply to the probabilities of hidden variables and tuples, and the joint token-latent tuple probability. For the latter, in particular, starting from the level- $(\ell + 1)$ hidden symbol $h_1^{(\ell+1)}$, lowest common ancestor (LCA) of X_1 and the tuple $\mathbf{h}_{s+1:2s}^{(\ell-1)}$, we have

$$\mathbb{P}\left\{X_{1}=\mu, \mathbf{h}_{s+1:2s}^{(\ell-1)}=\boldsymbol{\nu}\right\} = \sum_{\mu_{1},\dots,\mu_{\ell-1}=1}^{\nu} p_{1}^{(1)}(\mu|\mu_{1})\dots p_{1}^{(\ell)}(\mu_{\ell-1}|\mu_{\ell}) \times \sum_{\nu_{\ell-1},\mu_{\ell}} p^{(\ell)}(\boldsymbol{\nu}|\nu_{\ell}) p_{1,2}^{(\ell+1)}(\mu_{\ell},\nu_{\ell}|\mu_{\ell+1}) p_{1}^{(\ell+2)}(\mu_{\ell+1}).$$
(7)

For $\ell = 1$, the probability above coincides with the joint probability of the visible token X_1 and the tuple of visible tokens X_{s+1}, \ldots, X_{2s} . The correlations,

$$C^{(2)}(\mu, \boldsymbol{\nu}) \coloneqq \mathbb{P}\left\{X_1 = \mu, \mathbf{X}_{s+1:2s} = \boldsymbol{\nu}\right\} - \mathbb{P}\left\{X_1 = \mu\right\} \mathbb{P}\left\{\mathbf{X}_{s+1:2s} = \boldsymbol{\nu}\right\},\tag{8}$$

have been analyzed in Cagnetta and Wyart [6]: the mean vanishes, while the variance, in the limit of $m, v \to \infty$ with $f = m/v^{s-1}$ finite, follows

$$\left\langle \left(C^{(2)}(\mu, \boldsymbol{\nu}) \right)^2 \right\rangle \simeq \frac{(1-f)}{v^3 m^4}.$$
 (9)

For $\ell = 2$, after applying Equation 7, we get

$$C^{(3)}(\mu, \nu) = \sum_{\mu_1=1}^{v} p_1^{(1)}(\mu|\mu_1) \left(\mathbb{P}\left\{ h_1^{(1)} = \mu_1, \mathbf{h}_{s+1:2s}^{(\ell-1)} = \nu \right\} - \mathbb{P}\left\{ h_1^{(1)} = \mu_1 \right\} \mathbb{P}\left\{ \mathbf{h}_{s+1:2s}^{(\ell-1)} = \nu \right\} \right)$$
$$= \sum_{\mu_1=1}^{v} p_1^{(1)}(\mu|\mu_1) C^{(2)}(\mu_1, \nu), \tag{10}$$

where the last equality follows from noticing that the probability of level- ℓ hidden variables coincides with the probability of the leaves of a tree with $L - \ell$ layers. In general,

$$C^{(\ell+1)}(\mu, \boldsymbol{\nu}) = \sum_{\mu_1=1}^{\nu} p_1^{(1)}(\mu|\mu_1) C^{(\ell)}(\mu_1, \boldsymbol{\nu}),$$
(11)

thus

$$\left\langle \left(C^{(\ell+1)}(\mu, \boldsymbol{\nu}) \right)^2 \right\rangle = \sum_{\mu_1, \nu_1} \left\langle p_1^{(1)}(\mu | \mu_1) p_1^{(1)}(\mu | \nu_1) \right\rangle \left\langle C^{(\ell)}(\mu_1, \boldsymbol{\nu}) C^{(\ell)}(\nu_1, \boldsymbol{\nu}) \right\rangle$$
$$= \sum_{\mu_1} \left\langle \left(p_1^{(1)}(\mu | \mu_1) \right)^2 \right\rangle \left\langle \left(C^{(\ell)}(\mu_1, \boldsymbol{\nu}) \right)^2 \right\rangle + \sum_{\mu_1, \nu_1 \neq \mu_1} \left\langle p_1^{(1)}(\mu | \mu_1) p_1^{(1)}(\mu | \nu_1) \right\rangle \left\langle C^{(\ell)}(\mu_1, \boldsymbol{\nu}) C^{(\ell)}(\nu_1, \boldsymbol{\nu}) \right\rangle.$$
(12)

Knowing that the production rules of an RHM realization are chosen uniformly at random compatibly with the unambiguity constraint [6],

$$\left\langle \left(p^{(1)}(\mu|\mu_1) \right)^2 \right\rangle = \frac{v^{s-1}(v-1) + m(v^{s-1}-1)}{mv(v^s-1)},$$
(13)

and, for $\nu_1 \neq \mu_1$,

$$\left\langle p^{(1)}(\mu|\mu_1)p^{(1)}(\nu|\nu_1)\right\rangle = \frac{v^{s-1}-1}{v(v^s-1)}.$$
 (14)

In addition, since $\sum_{\mu} C^{(\ell)}(\mu, \nu) = 0$, then

$$\sum_{\nu_1 \neq \mu_1} \left\langle C^{(\ell)}(\mu_1, \boldsymbol{\nu}) C^{(\ell)}(\nu_1, \boldsymbol{\nu}) \right\rangle = -\left\langle \left(C^{(\ell)}(\mu_1, \boldsymbol{\nu}) \right)^2 \right\rangle.$$
(15)

Hence,

$$\left\langle \left(C^{(\ell+1)}(\mu, \boldsymbol{\nu}) \right)^2 \right\rangle = \frac{v^{s-1}(v-1)}{m(v^s-1)} \left\langle \left(C^{(\ell)}(\mu_1, \boldsymbol{\nu}) \right)^2 \right\rangle \xrightarrow{v \gg 1} \frac{1}{m} \left\langle \left(C^{(\ell)}(\mu_1, \boldsymbol{\nu}) \right)^2 \right\rangle.$$
(16)

Starting with $C^{(2)}$ from Equation 9, we get

$$C^{(\ell)} = \sqrt{\left\langle \left(C^{(\ell)}(\mu, \boldsymbol{\nu}) \right)^2 \right\rangle} \simeq \sqrt{\frac{(1-f)}{v^3 m^{2+\ell}}}.$$
(17)

Appendix D. One-step gradient descent

We consider a simplified one-step gradient descent setting [12], where a simple machine-learning model is trained to approximate the conditional probability of one input token X_{s+1} following an s-tuple of tokens $\mathbf{X} = (X_1, \ldots, X_s)$. The training set \mathcal{X}_P consists of P pairs (\mathbf{x}, ν) , with ν denoting the feature in the token X_{s+1} . We assume that

- *i*) the input tuple X is given as the one-hot encoding of the tuple index. Each of the mv possible combinations of s features is assigned an index $\mu = 1, ..., mv$ and x is the mv-dimensional sequence $x_{\mu} = \delta_{\mu,\mu(x)}$;
- *ii*) the machine-learning model is initialized on the empirical marginal probability of the token X_{s+1} over the training set, $\hat{\mathbb{P}}(X_{s+1} = \nu) := P^{-1} \sum_{(\boldsymbol{x},\lambda) \in \mathcal{X}_{\mathcal{P}}} \delta_{\nu,\lambda}$. This assumption is equivalent to a preprocessing step on the labels [12] that removes the class imbalance of the training set.

Due to assumption i), the task can be solved with a perceptron model followed by a softmax nonlinearity,

$$f_{\nu}(\boldsymbol{x};W) = \sum_{\boldsymbol{\mu}} W_{\nu,\boldsymbol{\mu}} \boldsymbol{x}_{\boldsymbol{\mu}}; \quad p_{\nu}(\boldsymbol{x};W) = e^{f_{\nu}(\boldsymbol{x};W)} \left(\sum_{\sigma} e^{f_{\sigma}(\boldsymbol{x};W)}\right)^{-1}; \quad (18)$$

where $W \in \mathbb{R}^{v \times (vm)}$ is the weight matrix. In this setup, Assumption *ii*) is realized by initializing the weights as $W_{\nu,\mu} = \log \hat{\mathbb{P}}(X_{s+1} = \nu)$ independently of μ .

The model f_{ν} of Equation 18 is trained via Gradient Descent on the empirical cross-entropy loss computed over a training set \mathcal{X}_P consisting of P pairs (\boldsymbol{x}, ν) , with ν denoting the feature in the token X_{s+1} ,

$$\mathcal{L} = \mathbb{E}_{(\boldsymbol{x},\nu)\in\mathcal{X}_P}\left[-\log\left(\frac{e^{f_{\nu}(\boldsymbol{x};W)}}{\sum_{\sigma=1}^{v}e^{f_{\sigma}(\boldsymbol{x};W)}}\right)\right],\tag{19}$$

where $\mathbb{E}_{(\boldsymbol{x},\nu)\in\mathcal{X}_P}$ denotes the empirical average over the training set. Denoting the learning rate with η , the update of the weights reads

$$\Delta W_{\nu,\mu} = -\eta \frac{\partial \mathcal{L}}{\partial f_{\nu}} \frac{\partial f_{\nu}}{\partial W_{\nu,\mu}} = \eta \mathbb{E}_{(\boldsymbol{x},\lambda)\in\mathcal{X}_{P}} \left[\delta_{\lambda,\nu} \boldsymbol{x}_{\mu} - \frac{e^{f_{\nu}}}{\sum_{\sigma=1}^{v} e^{f_{\sigma}}} \boldsymbol{x}_{\mu} \right]$$

$$= \eta \mathbb{E}_{(\boldsymbol{x},\lambda)\in\mathcal{X}_{P}} \left[\delta_{\lambda,\nu} \delta_{\mu,\mu(\boldsymbol{x})} - \hat{\mathbb{P}} \left(X_{s+1} = \nu \right) \delta_{\mu,\mu(\boldsymbol{x})} \right]$$

$$= \eta \left(\hat{\mathbb{P}} \left[X_{s+1} = \nu; \left(X_{1}, \dots, X_{s} \right) = (\mu_{1}, \dots, \mu_{s}) \right] - \hat{\mathbb{P}} \left[X_{s+1} = \nu \right] \hat{\mathbb{P}} \left[\left(X_{1}, \dots, X_{s} \right) = (\mu_{1}, \dots, \mu_{s}) \right] \right)$$

(20)

where, in the second line, we used assumption *i*) to replace x_{μ} with $\delta_{\mu,\mu(x)}$ and assumption *ii*) to replace $e^{f_{\nu}}/(\sum_{\sigma=1}^{v} e^{f_{\sigma}})$ with $\hat{\mathbb{P}}(X_{s+1} = \nu)$. The right-hand side of the last line equals the empirical token-tuple correlation $\hat{C}_{P}(\nu, \mu)$. Therefore, after one gradient step, the weights are given by

$$W_{\nu,\boldsymbol{\mu}} = \log \hat{\mathbb{P}} \left(X_{s+1} = \nu \right) + \eta \hat{C}_P(\nu, \boldsymbol{\mu}).$$
(21)

The first term is independent of the input μ , whereas the second can be thought of as a noisy measurement of the true token-tuple correlation $C(\nu, \mu)$. The true correlation is equal for all μ 's generated by the same higher-level hidden symbol $h^{(1)}(\mu)$ and its size can be estimated as the standard deviation over realizations of the RHM [6],

$$C^{(2)} = \left(\frac{1}{v^2 m} \frac{(1-f)}{v m^3}\right)^{1/2}.$$
(22)

The empirical measurement \hat{C}_P includes a sampling noise contribution, having size $(v^2mP)^{-1/2}$. If $P \gg P_2 = vm^3/(1 - f)$, then the \hat{C}_P in the right-hand side of Equation 21 is approximately equal to the true token-tuple correlation, thus the weights can be used to build a representation of the hidden variables of the generative model.

Appendix E. Experimental details

Random Hierarchy Model We train the U-Net-based Discrete Denoising Diffusion Probabilistic Model (D3PM), optimizing the diffusion loss derived from a variational bound on the negative log-likelihood [43]. Following Austin et al. [2], we use the neural network to predict the conditional expectation $\mathbb{E}(\boldsymbol{x}(0)|\boldsymbol{x}(t))$, which parameterizes the reverse diffusion process.

The convolutional U-Net consists of L resolution blocks in both the encoder and decoder, with a filter size of s, stride of s, and 8192 channels. Each block uses GeLU activation functions, and skip connections link encoder and decoder layers with the same resolution. The model also includes two embedding and unembedding layers, implemented as convolutions with filter size 1.

We initialize the network using the maximal-update (μ P) parameterization [48]. This allows stable feature learning dynamics even in large models. The model is trained with SGD with a learning rate of 1, using a batch size of 32, and momentum parameter of 0.9. The diffusion process follows a linear schedule with 1,000 noise levels. To prevent overfitting, we apply early stopping based on the validation loss, halting training when it plateaus or begins to increase.

Language diffusion model Our experiments are based on the codebase of MD4 [41]:

https://github.com/google-deepmind/md4. MD4 is a masked diffusion model. At each time step t, non-masked tokens either remain unchanged or transition to [MASK] with probability β_t . Using a one-hot-encoding representation of the $|\mathcal{V}| + 1$ states, the forward transition matrix is given by:

$$q(x_i(t)|x_i(t-1)) = (1-\beta_t)\mathbb{I} + \beta_t \mathbf{1}\mathbf{e}_M^{\top}, \tag{23}$$

with I the identity matrix, **1** a vector of ones and \mathbf{e}_M the one-hot-encoding vector corresponding to the [MASK] symbol. At the final time T, all tokens are masked, i.e., $x_i(T) = [MASK]$ for every $i \in [\dim(x)]$. We train MD4 with batch size 64 and context size 1024 on 4 H100s for a single epoch. All other hyperparameters are kept unchanged.

Vision diffusion model Our experiments are based on the codebase of Improved DDPMs [30]: https://github.com/openai/improved-diffusion. In particular, we train a DDPM with 128 channels, 3 resolution blocks, 4000 diffusion steps, cosine noise schedule, learning rate 10^{-4} and batch size 128 for 10 epochs using a *hybrid objective* [30].

Appendix F. Additional results

F.1. Dependence of sample complexity with m

To investigate the dependence of the accuracy on the number of synonyms m, we define the sample complexity P^* as the training set size at which the accuracy of the last level \mathcal{A}_L surpasses a threshold value \mathcal{A}^* . In our experiments, we set $\mathcal{A}^* = 1/2$.¹ Figure 4 shows the scaling behavior of P^* with m at fixed depth L = 2 (blue points). Empirically, we find good agreement with m^{L+1} (dashed line in the figure).

^{1.} Notice that the observed scaling of sample complexity remains robust to the specific choice of threshold value.



Figure 4: Sample complexity P^* for L = 2 in diffusion models and clustering algorithms based on correlations. Blue points show the empirical values of P^* for trained diffusion models, while black and red points represent clustering methods based on the correlations of latent tuples with the first token and the first visible tuple, respectively. The scaling $P^* \sim m^{L+1}$ aligns with theoretical predictions. Notably, the simple complexity of the diffusion model closely matches that of the correlation algorithm, suggesting that diffusion models learn hierarchical structures by leveraging statistical dependencies between synonyms.

F.2. Emergence of hierarchical representations in the U-Net

In Figure 5, we test the hypothesis that the U-Net learns to represent together inputs that differ by low-level synonyms, i.e., the choice of low-level production rules. To do so, we introduce a transformation operator $\mathcal{R}_{\ell} x$, which modifies a given data sample x by resetting all choices of the production rules emanating from layer ℓ . This operation is equivalent to substituting all tuples at depth $\ell - 1$ with a synonym. We then define the relative sensitivity $S_{k,\ell}$ of the pre-activations a_k at layer k to the transformation \mathcal{R}_{ℓ} :

$$S_{k,\ell} = \frac{\mathbb{E}_{\boldsymbol{x}}[\|a_k(\boldsymbol{x}) - a_k(\mathcal{R}_{\ell}\,\boldsymbol{x})\|^2]}{\mathbb{E}_{\boldsymbol{x},\boldsymbol{y}}[\|a_k(\boldsymbol{x}) - a_k(\boldsymbol{y})\|^2]}.$$
(24)

Here, the numerator measures how much the activations change when synonym substitutions are applied at depth ℓ , while the denominator normalizes by the overall variability of activations across different data points. A low value of $S_{k,\ell}$ indicates that the network is invariant to synonym substitutions at depth ℓ , implying that it has learned the corresponding compositional rule.

Figure 5 shows the relative sensitivity of each layer as a function of the number of training points P. As P increases, the sensitivities $S_{k,\ell}$ decrease sequentially across levels, following the same staged learning process observed in Figure 1. Deep encoder layers become invariant to synonym substitutions at lower levels, confirming that the network is learning to encode the hierarchical structure of the grammar. In contrast, decoder layers gradually regain sensitivity to specific low-level symbols as the output is approached. This behavior aligns with their role in reconstructing low-level details from high-level representations. Crucially, the network begins to satisfy rules at level ℓ precisely when it becomes insensitive to synonymic variations at level $\ell - 1$. This suggests



Figure 5: Relative sensitivity of the hidden representations of the U-Net, defined in Equation 24, with respect to the number of training points P. Different colors correspond to different levels ℓ of synonymic exchange, while different panels correspond to the pre-activations of different U-Net blocks. Encoder layer 1 is the closest to the input, while decoder layer 5 is the closest to the output. As the number of training points increases, deeper layers of the encoder become less sensitive to deeper synonymic transformations. This implies that deeper encoder layers learn to represent deeper latent variables of the RHM. The decoder layers, instead, progressively regain the sensitivity to the synonyms layer-by-layer as they expand latent variables into their lower-level representations. For each level ℓ , the dashed line represents the fraction of generated samples that do not satisfy the rules at that level, i.e., $1 - A_{\ell}$. The U-Net learns to satisfy rules at level ℓ when it becomes insensitive to the synonyms of the variables at level $\ell - 1$.



Figure 6: Sample complexity of clustering with L = 3. Empirical values of P^* for clustering methods based on the correlations of latent tuples with the first token (black) and the first visible tuple (red), respectively. The scaling $P^* \sim m^{L+1}$ aligns with theoretical predictions.

that the U-Net learns to collapse lower-level synonyms into shared latent representations and to compose these latents according to the production rules at level ℓ .

F.3. Sample complexity of deep clustering algorithm

In Figure 6, we test our theoretical prediction for the hierarchical clustering algorithm with L = 3. Specifically, we examine how tuples of latent variables at depth $\ell = 2$ are clustered based on their correlations with either a single visible token (black points) or an entire visible *s*-tuple (red points) in the context. As predicted in Section B, the sample complexity of both clustering approaches scales as m^4 , confirming our theoretical result.

F.4. Perplexity of the generated text

Figure 7 presents an alternative measure to correlations in the generated text for quantifying the longer and longer coherence as training progresses. Specifically, we extract sentences from the generated datasets and estimate token-level average log-likelihoods using LLaMA-2-7B [46], i.e., we compute

$$\mathbb{E}_{x_{0:T}}[\log p_{\text{LLM}}(x_T | x_{0:T-1})]$$
(25)

for a token x_T as a function of its context length T. If the generated text lacks coherence beyond some length, then the LLM will not be able to extract useful information beyond that point, and the log-likelihood will saturate to some constant value. Figure 7 reports the corresponding *perplexity*, defined as the exponential of the negative log-likelihood (25), where the average is done over 1024 samples. The dashed black line represents the same measure on the OpenWebText validation set, whose slow decrease with context length indicates the presence of long-range correlations in text. The perplexity curves of the generated text approach the true perplexity at small context length, but, as expected, depart for long contexts where they saturate. Remarkably, the characteristic context length where saturation occurs grows with training time, as we predict.



Figure 7: **Perplexity of the generated text as a function of the conditioning context length computed with LLaMA-2-7B.** Averages done over 1024 samples. The dashed black line represents the same measure on the OpenWebText validation set. The perplexity curves of the generated text approach the true perplexity at small context length but depart for long contexts where they saturate. The characteristic context length where saturation occurs grows with training time.

F.5. Depth-wise maximum mean discrepancy

To quantify hierarchical learning of vision diffusion models, we analyze the hierarchical and compositional structure of generated images using deep latent representations from a pre-trained ResNet-18 [18]. Early layers encode low-level localized features, while deep layers represent more abstract and global factors [25, 33], as also observed for CNNs trained on the RHM [8]. We compute the *Maximum Mean Discrepancy* (MMD) [17] between ResNet embeddings of the generated images and those from the ImageNet validation set. MMD-based evaluations with deep network embeddings have recently been proposed as a robust metric for assessing image quality in diffusion models [21].

Figure 8b presents the MMD measured at different depths of the ResNet model as a function of the number of seen examples. Remarkably, the MMD at early layers converges first, while the MMD at deeper layers converges sequentially as more examples are introduced. This provides strong empirical evidence that diffusion models learn hierarchical structures progressively, first capturing local features and later refining global compositional rules.



(a) Images generated at different training stages.



Figure 8: **Stage-wise learning of vision diffusion model on ImageNet64.** (a) Examples of images generated by the diffusion model at different training steps. (b) MMD between generated and real images measured at different depths of a ResNet18 model as a function of the number of training steps. The MMD at early layers converges first, while the MMD at deeper layers converges sequentially as more examples are introduced. The grey dashed line indicates the end of the first epoch.